# Final Project Report for IAL620-02: Text Mining & Natural Language Processing

# Water Reports & News Classification using BERT Classification and sentiment analysis, Text mining, word embedding.

Pradnya Patil

Master's in informatics and Analytics
University of North Carolina, Greensboro, USA
papatil@uncg.edu

## Problem Statement:

Now a days due to global warming and high amount of pollution, we are getting new information about environment and that is accessible to public on every second, the information such as reports, news, books, articles, etc. Using the water report and news, we can classify the data into different category and make the decision more efficiently.

In this project, we seek to answer the following data research questions:

A. Exploratory Data Analysis on the raw downloaded dataset along with text cleaning and transformation.

B. Explain different methods to analyze text and extract features that can be used to build a classification model

C. Classify water news and Reports based on subject text using Machine Learning Models (Decision Trees, SVM, Random Forest...).

D. Predict the probability of the different model of classification.

E. Do sentiment analysis of news and reports related to the water problems based on environmental issue.

## Project Abstract:

In last years, water quality has been threatened by various pollutants and some environmental issues. Therefore, modeling and predicting water reports and news have become very important in controlling water pollution. In this project, models are built to classify the reports and news based on the 5 different target variables using the machine learning algorithms.

For this project, I am collecting the from https://www.kaggle.com/vbmokin/nlp-reports-news-classification. From this analysis, I will learn the key concept of the NLP mostly TEXT classification where the problem of assigning the news and report into their categories. Along with the Classification, I will learn the sentiment analysis also.

From this process, we can tell the reports, or the news are related to environment problems, pollution, treatment, climate indicators or the biotic monitoring in water or in a river basin.

## Data Source and description:

The dataset contains 2 files so far - an English-language dataset from the English-language edition of the book, where the co-author, and a Ukrainian-language dataset from a separate Ukrainian-language edition of this book. These datasets contain approximately 95% of the same information.

Text Feature:
   • Text - One or more sentences from reports or news

Binary Target Features:
   • Env_problems - Is the text about an environmental problem? (0 or 1)
   • Pollution - Is the text about environmental pollution? (0 or 1)
   • Treatment - Is the text about treatment plants or environmental technologies? (0 or 1)
   • Climate - Is the text about climatic indicators? (0 or 1)
   • Biomonitoring - Is the text about biological, biotic monitoring in water or in a river basin? (0 or 1)

In this project I have used the following languages, frameworks, tools, libraries, packages for data extraction to Sentiment analysis and water reports & news data classification/regression.

**Technologies:** *R and Python*

 **Libraries:**

R - tidytext, tidyverse, tm, lubridate, ldatuning, fliptime, stringr, ggmap, maps, stringr, rvest, quantenda, readtxt, textplots, widyr, ggplot2, igraph, ggraph

Python - numpy, pandas, matplotlib, seaborn, sklearn, torch, transformers, summarizer

**Tools:**  Rstudio, Jupyter hub for python

**Methods:**  R Selenium, BERT Classification

**Water Report & News Data Description Table -**

| Attribute | Data Type | Description |
|---|---|---|
| text | String | One or more sentences from reports or news |
| env_problems | Integer | Is the text about an environmental problem? (0 or 1) |
| pollution | Integer | Is the text about environmental pollution? (0 or 1) |
| treatment | Integer | Is the text about treatment plants or environmental technologies? (0 or 1) |
| climate | Integer | Is the text about climatic indicators? (0 or 1) |
| biomonitoring | Integer | Is the text about biological, biotic monitoring in water or in a river basin? (0 or 1) |

**WSJ News Data Description Table -**

| Attribute | Data Type | Description |
|---|---|---|
| URL | String | The URL path of the news site |
| Title | String | The Title of the news from that URL |
| Date | String | The date of the news that published on what day, and by whom. |
| Text | String | All the Text contains the news that published by WSJ. |

## Data Analysis & Results –

News, reports, and articles are the most important unstructured data sources now a days along with the structured dataset. In the news, reports people describe about the health, food, world, environment, travel, and every topic. From this unstructured data form, we can in out the different result and do the data analysis and help to the world in different aspects.
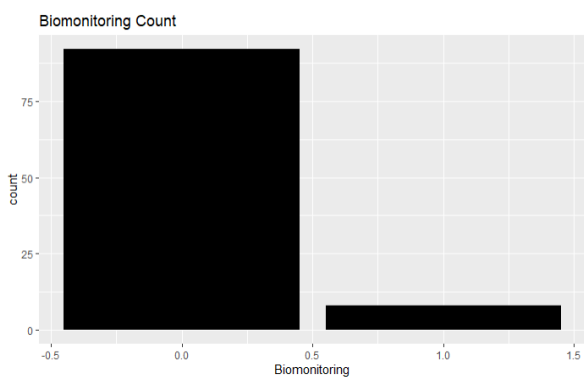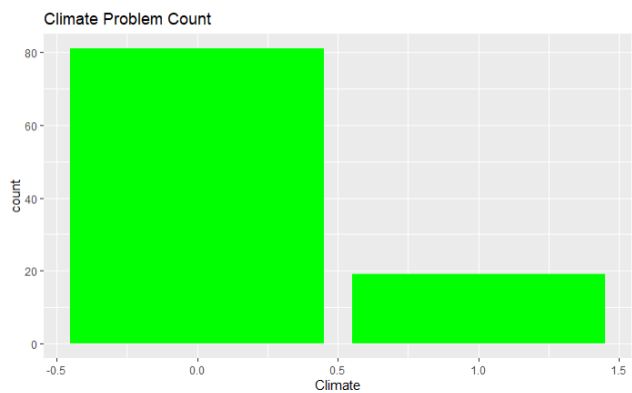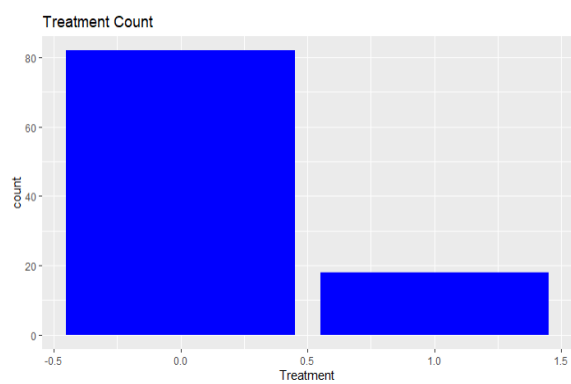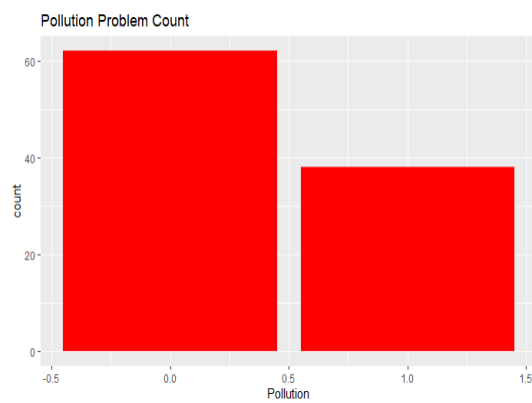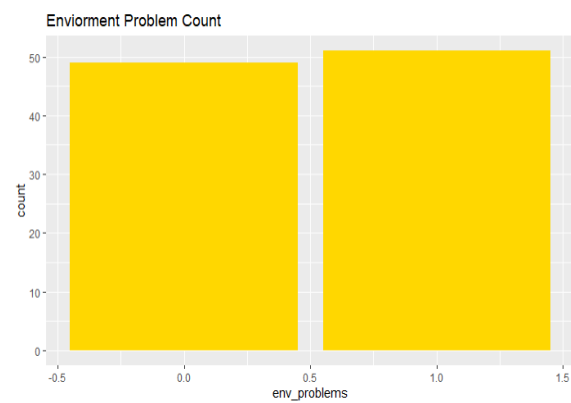
Here are some of my analyses about the news and reports from my Exploratory data analysis.
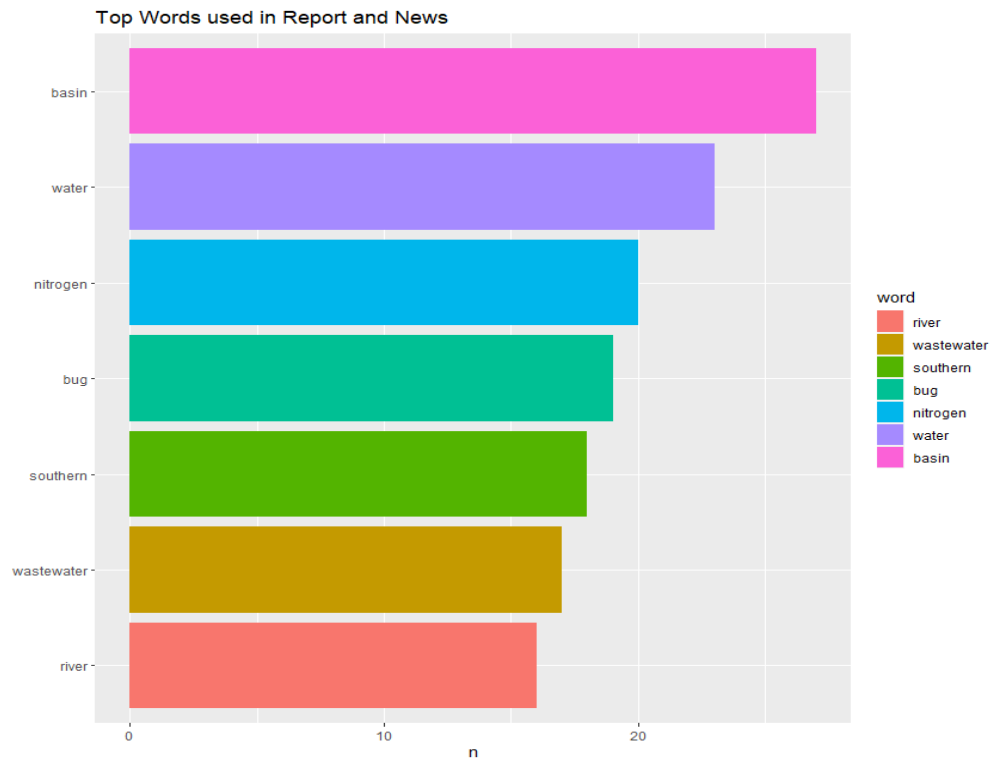
**Water Report and News EDA -**

This is the graph that shows the missing values of your data set. In my Water Report and News data set I don't have any missing values.

Figure 1

The below graphs are the bar plot of the target variables from the dataset.

This is the diagram of the top words that have been used in the text of the dataset.



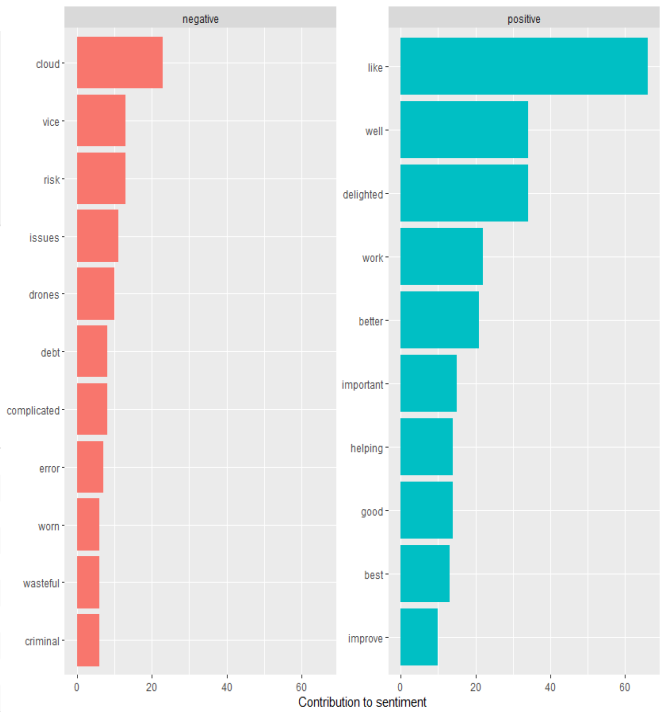Top Words used in Report and News



```
# capture negative sentiments for bing lexicon
# chapter 2.2 Text Mining with R
bing.negative <- get_sentiments("bing") %>%
  filter(sentiment == "negative")

bing.negative
```
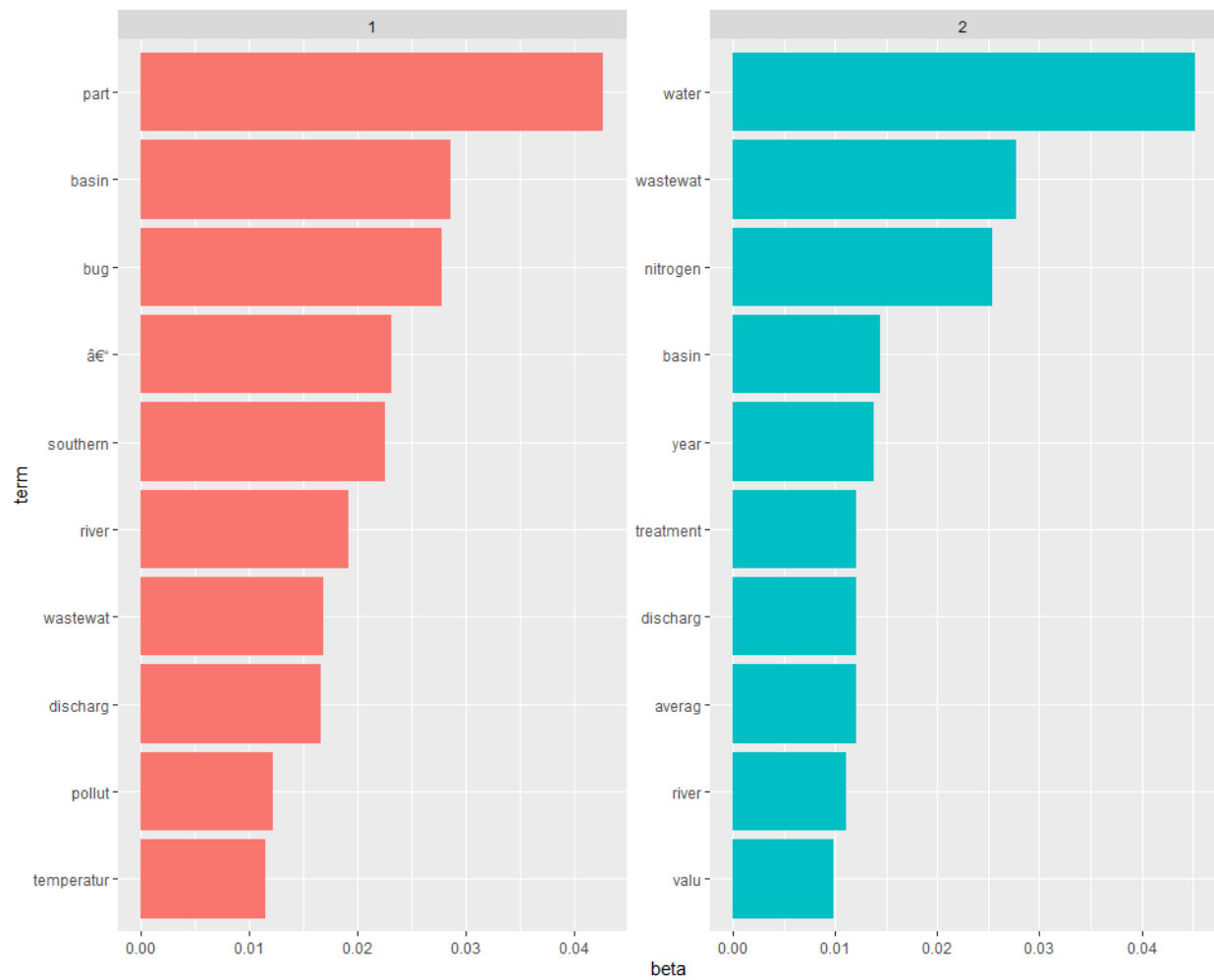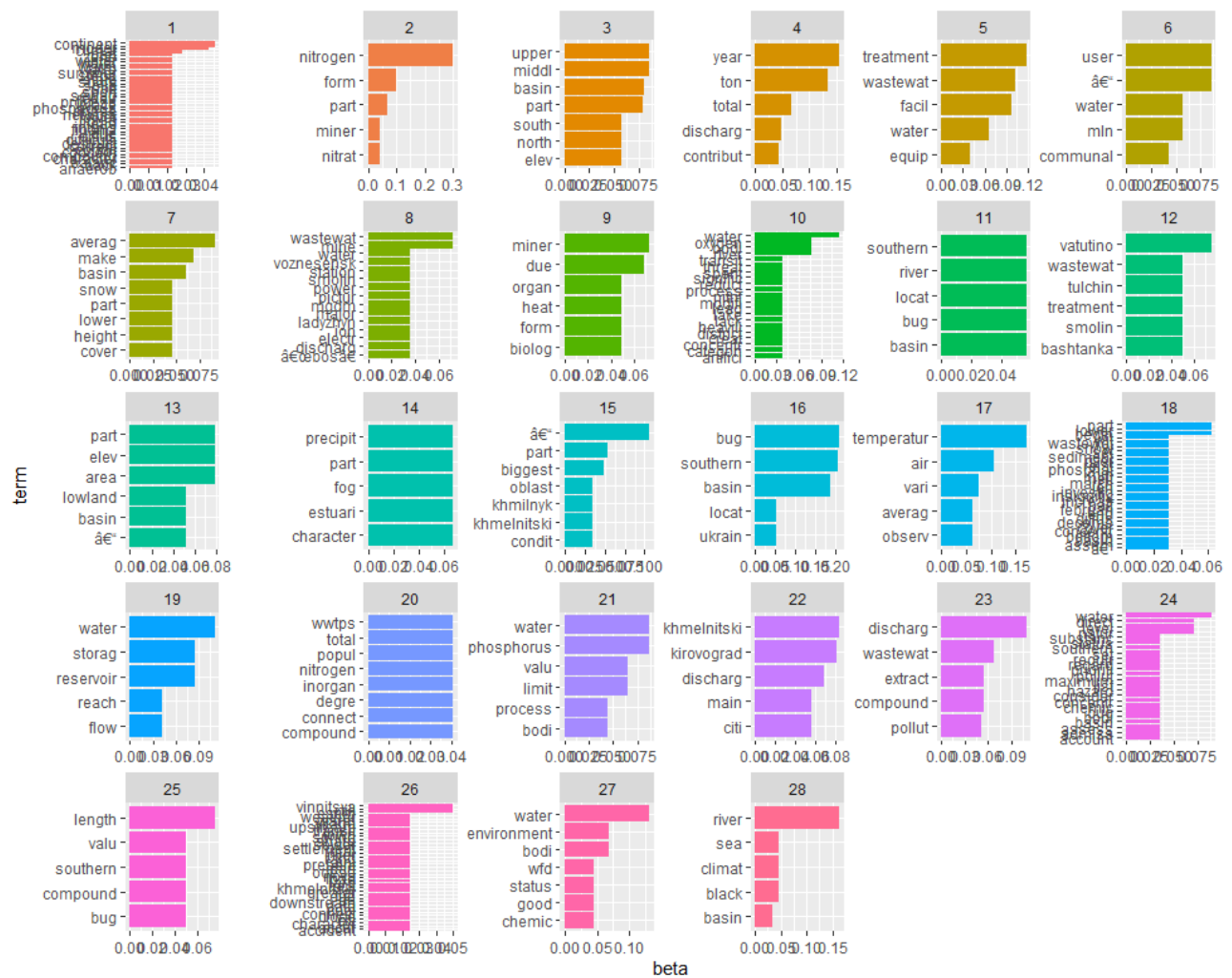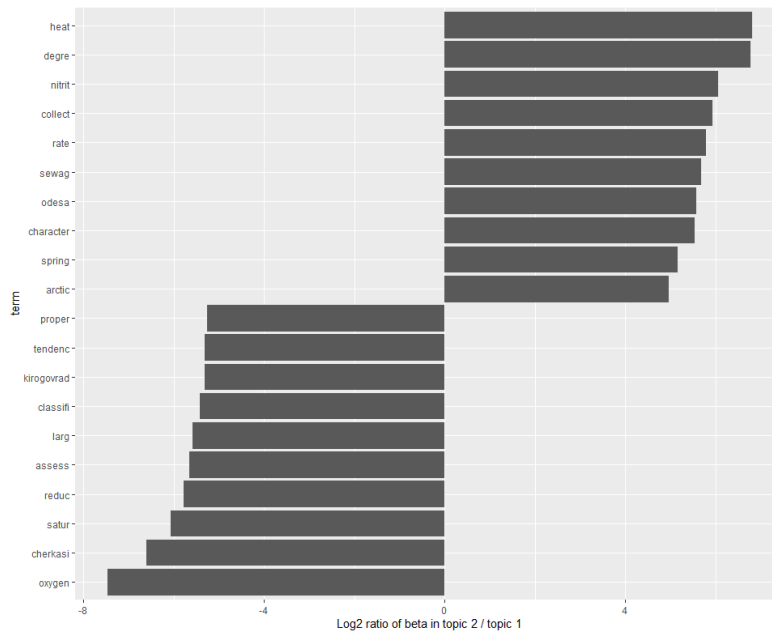```
A tibble: 4,781 x 2
```

| word <chr> | sentiment <chr> |
|---|---|
| 2-faces | negative |
| abnormal | negative |
| abolish | negative |
| abominable | negative |
| abominably | negative |
| abominate | negative |
| abomination | negative |
| abort | negative |
| aborted | negative |
| aborts | negative |

A tibble: 1,035 x 2

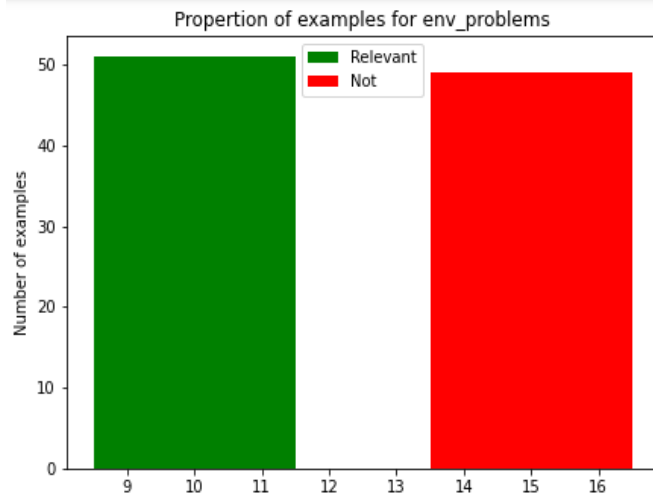| bigram<br><chr> | n<br><int> |
|---|---|
| southern bug | 17 |
| bug basin | 11 |
| water bodies | 9 |
| lower part | 6 |
| air temperature | 5 |
| surface waters | 5 |
| â mln | 4 |
| basin located | 4 |
| khmelnitsky kirovograd | 4 |
| mln users | 4 |

1-10 of 1,035 rows

## BERT Classification Results and Plots -


Propertion of examples for env_problems
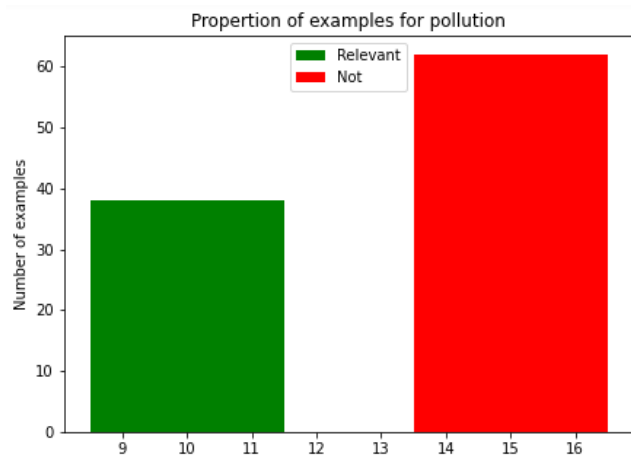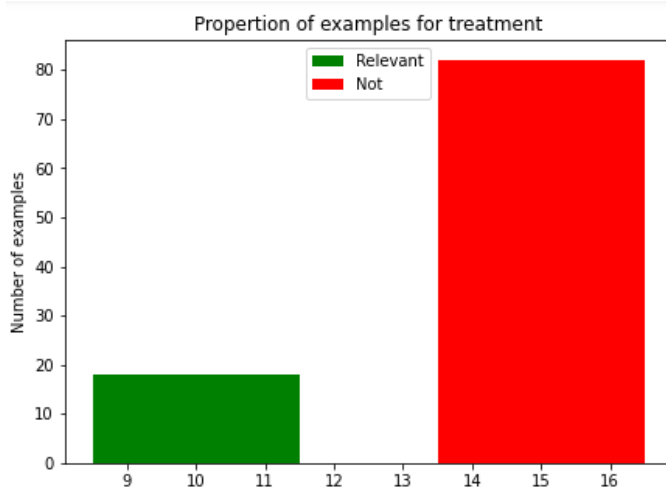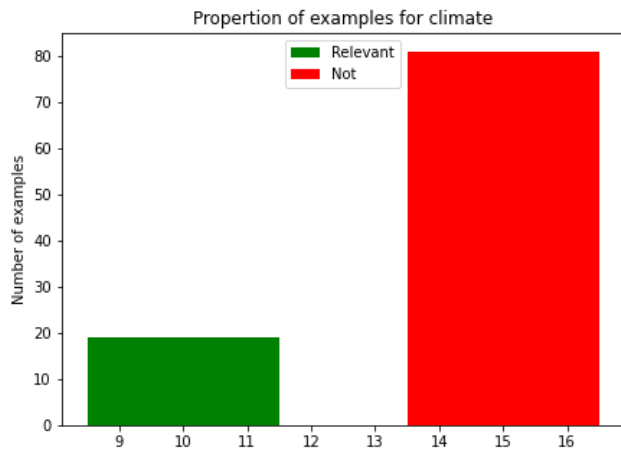
Classification for env_problems:
best parameters: {'C': 5.263252631578947}
best scores: 0.8
Score of the test prediction - 0.9


Propertion of examples for pollution
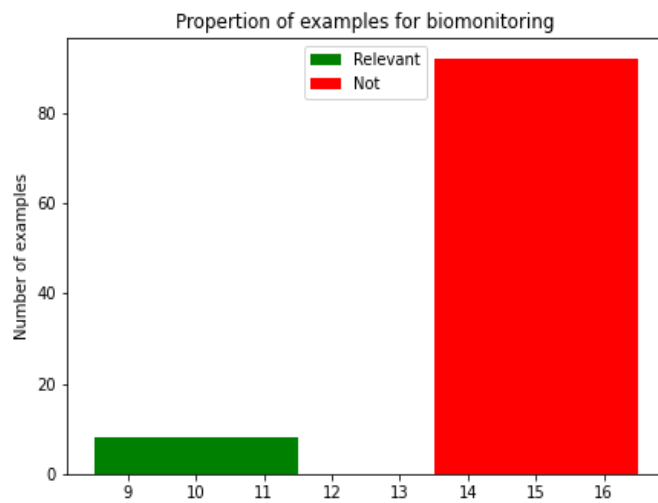
Classification for pollution:
best parameters: {'C': 5.263252631578947}
best scores: 0.8166666666666668
Score of the test prediction - 0.7


Propertion of examples for treatment

Classification for treatment:
best parameters: {'C': 0.0001}
best scores: 0.8166666666666668
Score of the test prediction - 0.825

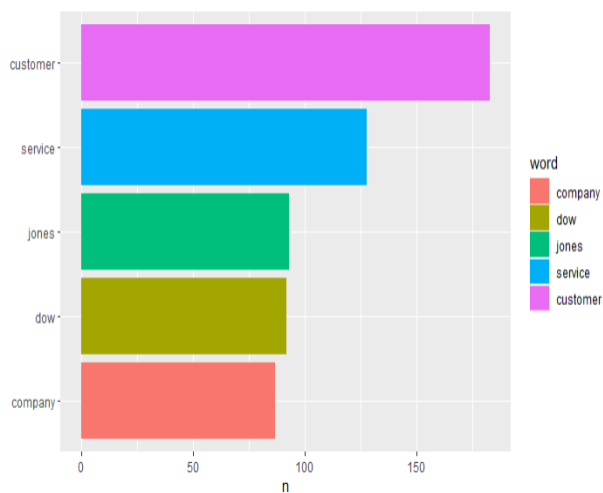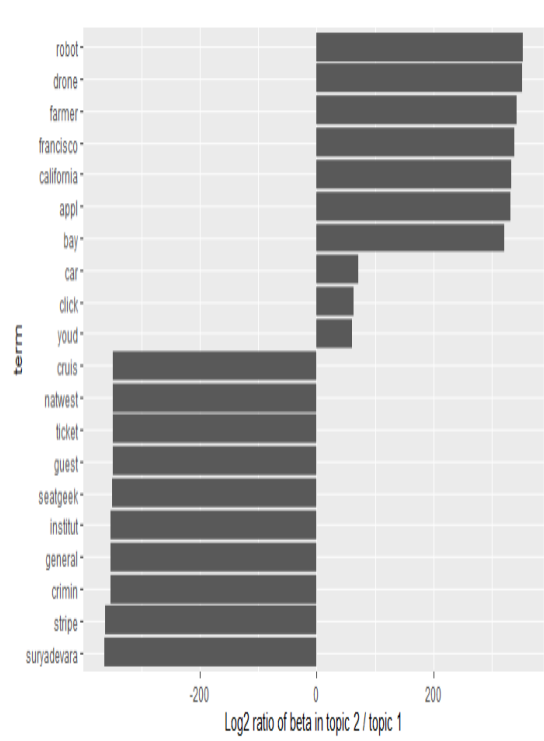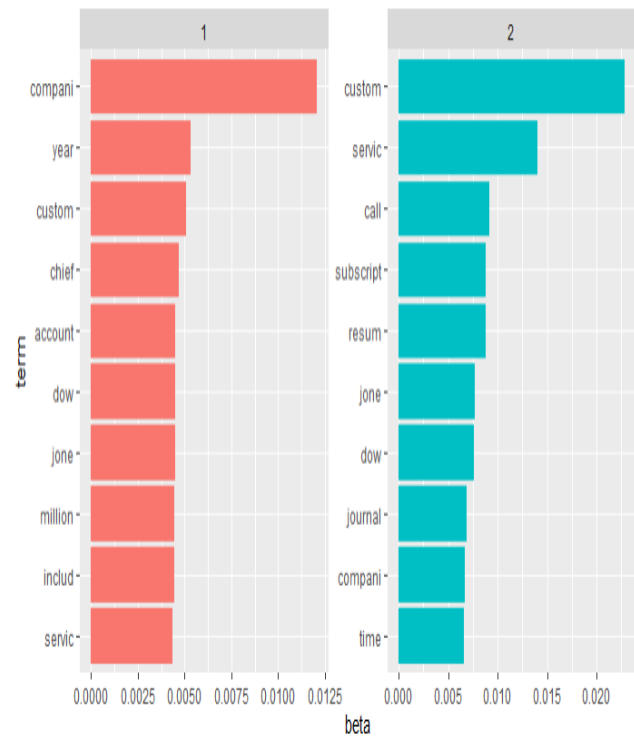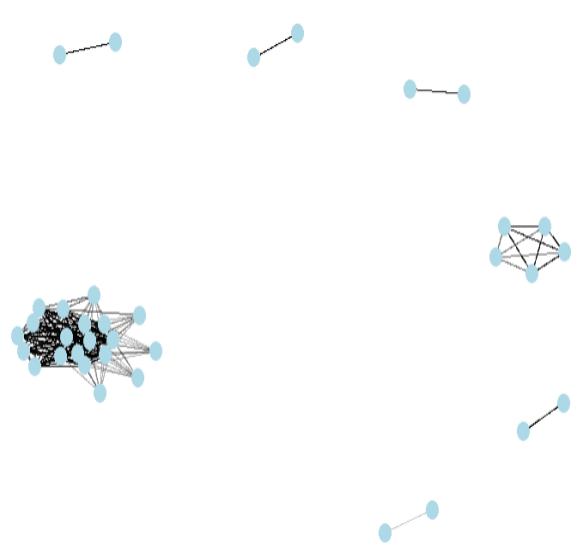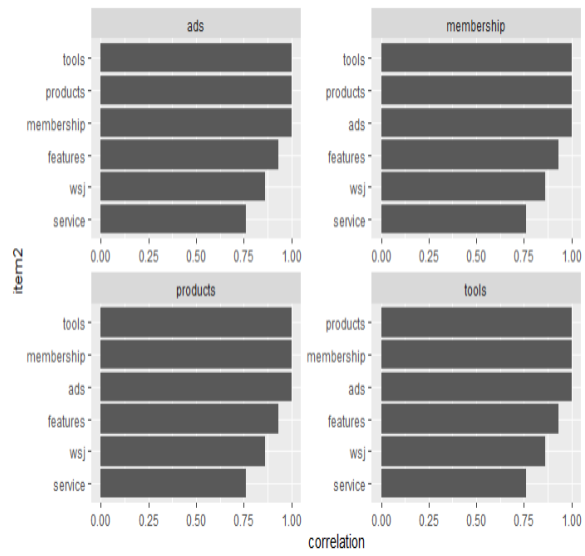Proportion of examples for climate

Classification for climate:
best parameters: {'C': 5.263252631578947}
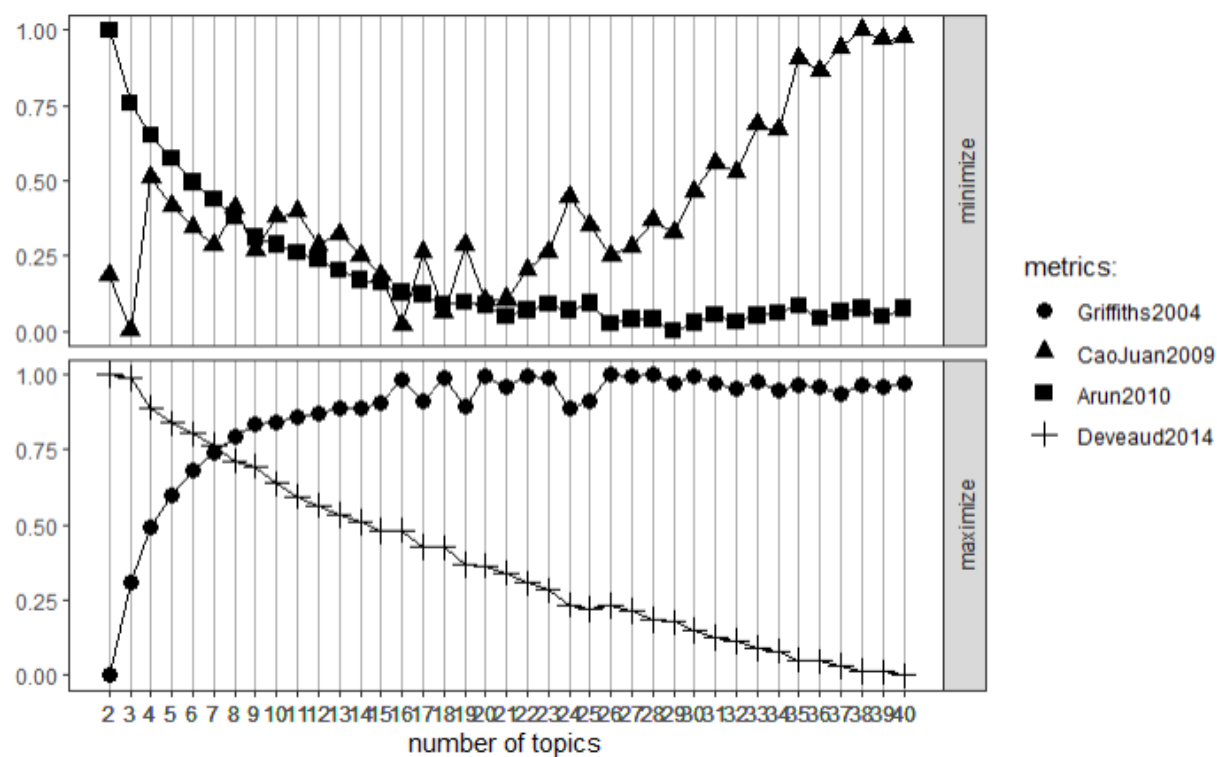best scores: 0.9666666666666666
Score of the test prediction - 0.975



Proportion of examples for biomonitoring

Classification for biomonitoring:
best parameters: {'C': 0.0001}
best scores: 0.95
Score of the test prediction - 0.875

## WSJ News EDA –

**negative**

| | Contribution to sentiment |
|---|---|
| cloud | |
| vice | |
| risk | |
| issues | |
| drones | |
| debt | |
| complicated | |
| error | |
| worn | |
| wasteful | |
| criminal | |

**positive**

| | |
|---|---|
| like | |
| well | |
| delighted | |
| work | |
| better | |
| important | |
| helping | |
| good | |
| best | |
| improve | |

Words preceded by "not"

| | Sentiment value * number of occurrences |
|---|---|
| responsible | |
| like | |
| accept | |
| disappear | |
| afraid | |

**no**

Words preceded by negation term

| | |
|---|---|
| benefits | |
| clear | |

**not**

| | Sentiment value * # of occurrences |
|---|---|
| responsible | |
| like | |
| accept | |
| disappear | |
| afraid | |

**without**

| | |
|---|---|
| pay | |

## Conclusion:

Github URL for code: https://github.com/Pradapatil/IAL62_NLP_-_TM

From this subject and Project, I have learned Text Mining and the Natural language processing methods. I did some data wrangling and scrubbing and find out the most frequent words that has been used in the reports and news.

From the Sentiment analysis learned how negative and positive words affected on our analysis. Using the word frequencies, N-gram an word correlation understood the concept and the relation between the words that has been used in text format.

Using the BERT classification model got the best score (0.96) for the climate variable and the got the highest prediction for the same variable which is  (0.975)

## References:

- Sklearn, Python library

**References:**

• https://www.kaggle.com/vbmokin/nlp-reports-news-classification

• https://towardsdatascience.com/text-analysis-feature-engineering-with-nlp-502d6ea9225d